

Introducción a la pangenómica microbiana

usando GET_HOMOLOGUES

Pablo Vinuesa

Centro de Ciencias Genómicas – UNAM,
vinuesa[at]ccg.unam.mx

<http://www.ccg.unam.mx/~vinuesa/>

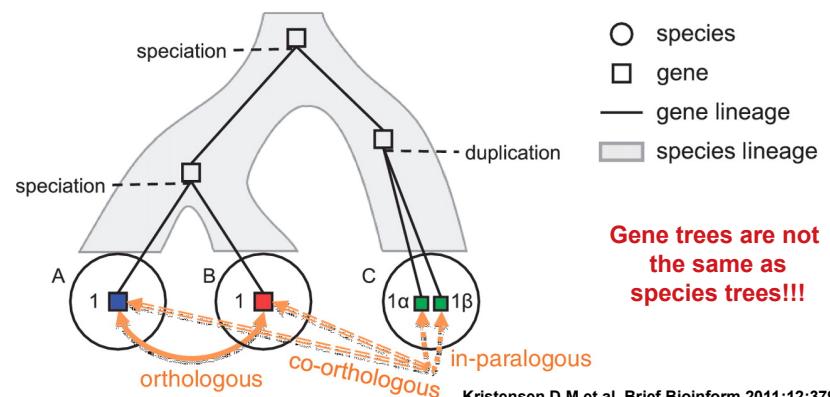
Introducción a la filogenómica y pangenómica microbiana,
TIB2019-T3, 1 Marzo 2019

<https://github.com/vinuesa/TIB-filoinfo>



Computational methods to identify homologous gene families

Orthology, co-orthology and paralogy relationships in the evolution of four genes that arose from a single common ancestor.



Orthologs: essentially instances of 'the same gene' in different species. Tend to retain function across large phylogenetic distances. Key markers for phylogenetics and genome annotation.

Paralogs: Tend to diverge over time to perform different functions via subfunctionalization or neofunctionalization routes.

Outline

I. Basic concepts in microbial comparative (pan-)genomics

- Methods for orthology inference
- Genome mosaicism, genomic islands and HGT
- The prokaryotic cloud, shell and core genes
- Microbial core- and pan-genomes

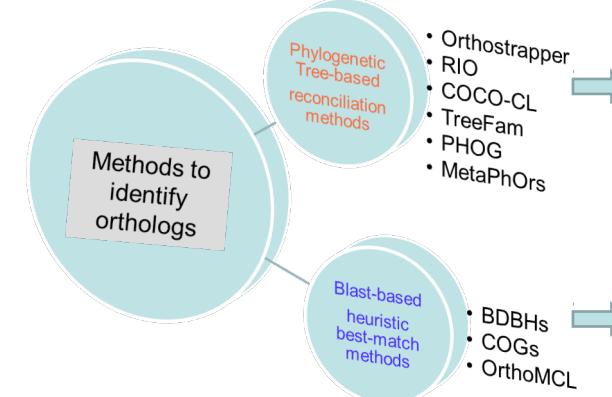
II. get_homologues, a powerful and highly configurable software package for microbial pan-genomics

- Overview of the package's capabilities (get_homologues.pl and accessory scripts)
- Using GET_HOMOLOGUES to explore factors affecting the calling of homologues

III. A pangenomic analysis of pIncA/C plasmids using GET_HOMOLOGUES

- Defining a robust core- and pan-genome of pIncA/C plasmids
- Exploring the gene space of pIncA/C plasmids: core, shell, cloud
- Pan-genome trees vs. core genome trees (supermatrices)
- Identifying lineage-specific genes in NDM-1 producing plasmids

Computational methods to identify homologs



- Most direct approach but known to be subjected to increasing artifacts as evol. dist. increases due to:
 - alignment problems
 - long-branch attraction
 - heterotachy
 - tree searches stuck in local sub-optimal maxima
- More adequate for bacteria, as they don't follow a clear tree-like evolutionary path (HGT).
- The only practical strategy for large datasets, therefore most common method in comparative genomics

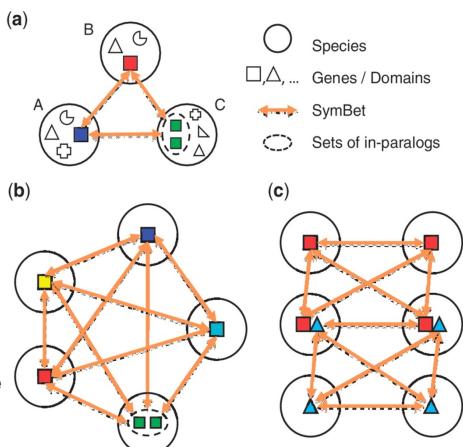
Computational methods to identify orthologs: BBHs

Grouping of genes in different species that are each others' BLAST BBHs into sets of orthologs and co-orthologs.

Linking pairs of BBHs from multiple genomes has a property of **self-verification**, as their consistency would be very unlikely due to chance, especially between phylogenetically distant lineages.

Methods for clustering of pairwise BBHs vary, but the most widely used approach involves a **single-linkage clustering procedure**, where any two clusters sharing a common BBH are merged until convergence.

Problems with: differential gene loss, domain recombination/gain/loss



Kristensen D M et al. Brief Bioinform 2011;12:379-391

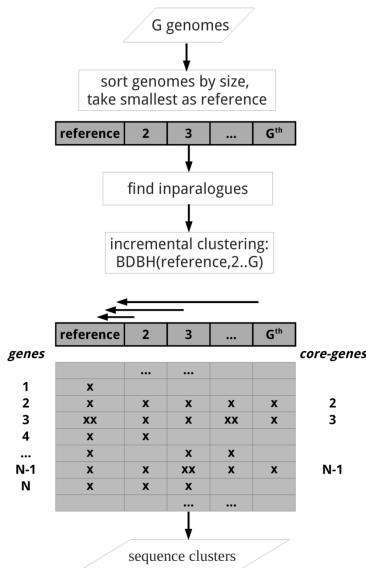
Computational methods to identify orthologs: BDBHs

The bidirectional best-hit approach (BDBH)

As implemented in GET_HOMOLOGUES

(Contreras-Moreira & Vinuesa 2013)

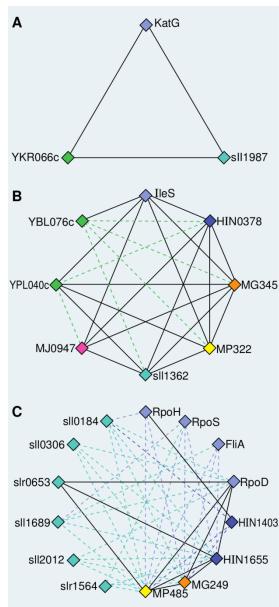
- Find inparalogues in all genomes
- Find reciprocal best blast hits between reference and remaining genomes
- Cluster sequences based on filtering criteria such as:
 - % alignment coverage
 - min. % sequence identity
 - E-score cut-off
 - Pfam domain-composition
 - synteny



Computational methods to identify orthologs: COGtriangles

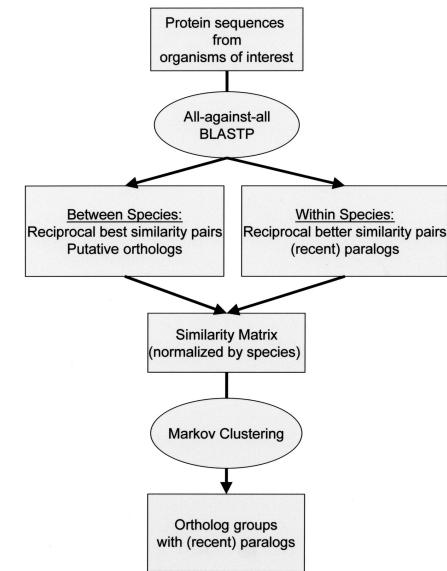
Method	Description
Clusters of orthologous groups (COGs), variants and derivatives	Identifies three-way BBHs between orthologs or sets of co-orthologs in three different species , and these groups expanded (merging triangles whenever they share a common side) until saturation, followed by manual splitting of large groups improperly joined by multidomain proteins or complex mixtures of in- and out-paralogs.

Tatusov et al. Science 1997. Vol. 278:631-637
Kristensen et al. Bioinformatics 2010. Vol. 26:1481-1487



Computational methods to identify orthologs: OrthoMCL

Method	Description
OrthoMCL	Forms groups of orthologs and co-orthologs using a Markov clustering process involving iterative simulations of stochastic (randomized) flow on the edges of a BBH graph , with clusters of desired tightness identified depending on a given 'inflation' parameter determined by trial and error.



Li L et al. Genome Res. 2003;13:2178-2189

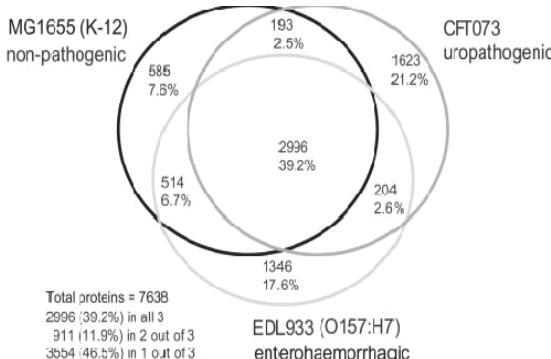
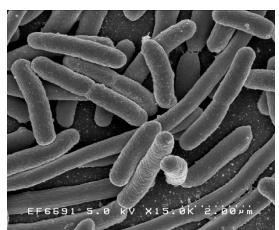
Basic concepts in bacterial comparative genomics - Mosaic genome structure and the pan-genome

Extensive mosaic structure revealed by the complete genome sequence of uropathogenic *Escherichia coli*

R. A. Welch*, V. Burland†‡, G. Plunkett III‡, P. Redford*, P. Roesch*, D. Rasko§, E. L. Buckles§, S.-R. Liou†, A. Boutin†**, J. Hackett††, D. Stroud*, G. F. Mayhew*, D. J. Rose*, S. Zhou††, D. C. Schwartz§§, N. T. Perna§§, H. L. T. Mobley§, M. S. Donnenberg*, and F. R. Blattner*

17020–17024 | PNAS | December 24, 2002 | vol. 99 | no. 26

Edited by John J. Mekalanos, Harvard Medical School, Boston, MA, and approved October 24, 2002 (received for review August 30, 2002)



Basic concepts in bacterial comparative genomics - Mosaic genome structure and the pan-genome

"accessory/flexible genome (shell + cloud genes)":

- patchy distribution among strains
- ecological specialization
- Plasmids, phages, GIs, ICE, integrons, ISs ...
- highly transferable (HGT)

"core genome":

- orthologous gene families
- species/population trees
- population genetics
- Genome annotation

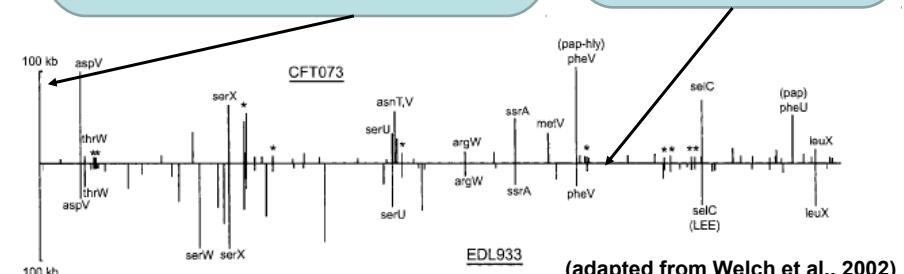
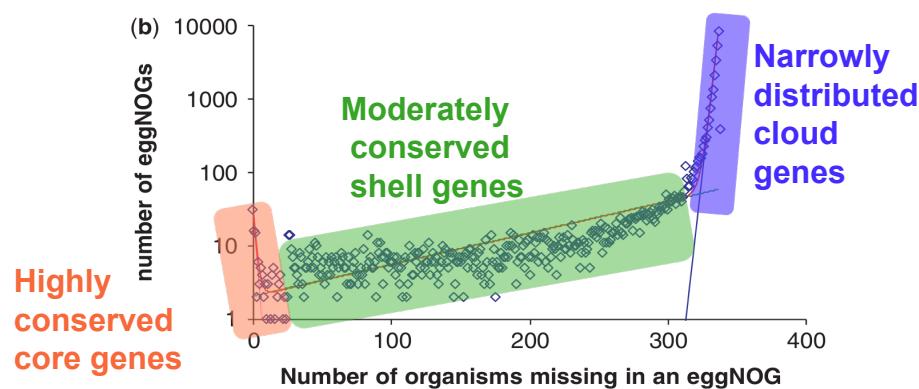


Fig. 3. Locations and sizes of CFT073 and EDL933 islands. Island size, vertical axis; position in colinear backbone, horizontal axis. All islands >4 kb are shown. Islands located at tRNAs are indicated by tRNA labels. One tmRNA (ssrA) is also an insertion target. *, CFT073 and EDL933 islands in the same backbone location but not near tRNAs.

Basic concepts in bacterial comparative genomics - Mosaic genome structure and the pan-genome

Distribution of COGs by the number of organisms included in each cluster:



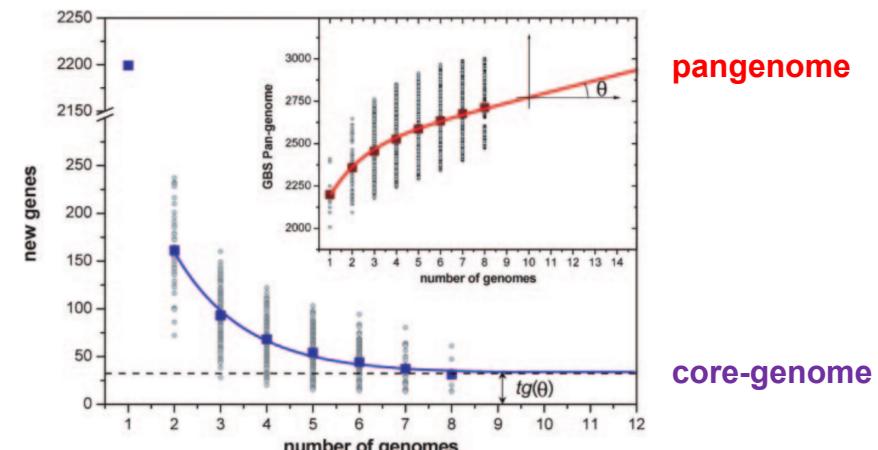
Fit of an exponential decay function with three exponents:

The core (~70 clusters), shell (~5700 clusters) the cloud (~ 24,000 clusters)

Koonin & Wolf, NAR 2008. Vol 36, No. 21

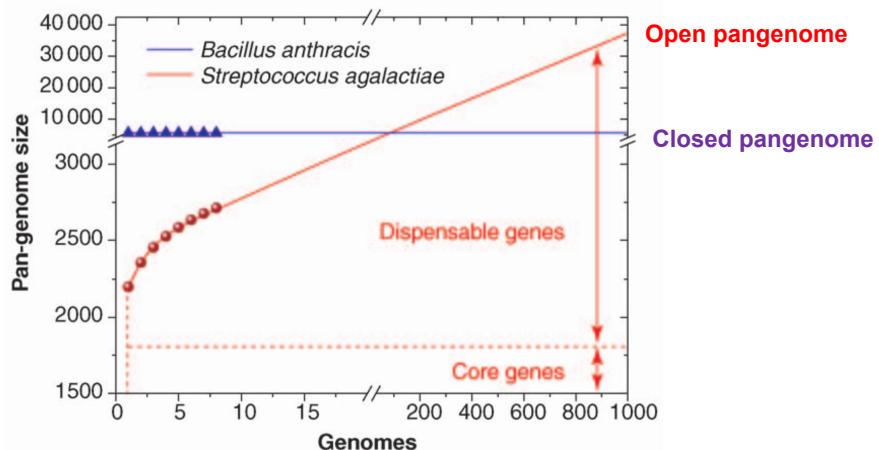
Basic concepts in bacterial comparative genomics - the pan-genome (Tettelin et al. 2005)

The core and pan-genome for 8 *Streptococcus agalactiae* genomes



Tettelin, H. et al. (2005). *Proc. Natl. Acad. Sci. U. S. A.* **102**(39): 13950-13955.

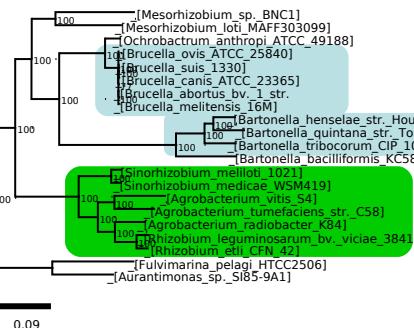
Basic concepts in bacterial comparative genomics - the pan-genome (Tettelin et al. 2005)



Tettelin, H., et al. (2008). "Curr. Opin. Microbiol. 11(5): 472-7.

Basic concepts in bacterial comparative genomics - the pangenome (Tettelin et al. 2005)

Size estimations of the core and pangenomes of the animal-associated Rhizobiaceae (5 *Brucella* spp., 4 *Bartonella* spp.)



Size estimations of the core and pangenomes of the plant-associated Rhizobiaceae (2 *Rhizobium*, 2 *Sinorhizobium*, 3 *Agrobacterium*)

Vinuesa & Contreras-Moreira, in prep.

GET_HOMOLOGUES, a Versatile Software Package for Scalable and Robust Microbial Pangenome Analysis



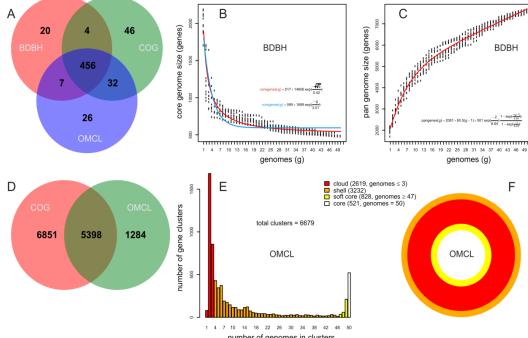
Bruno Contreras-Moreira,^{a,b} Pablo Vinuesa^c

Estación Experimental de Aula Dei (EEAD-CSIC), Zaragoza, Spain^a; Fundación ARAID, Zaragoza, Spain^b; Centro de Ciencias Genómicas, Universidad Nacional Autónoma de México, Cuernavaca, Morelos, Mexico^c



Applied and Environmental Microbiology p. 7696–7701 December 2013 Volume 79 Number 24

1. Clustering of prot. and nt secs. (CDSs) in homologous groups (paralogs, orthologs).
2. Identification of orthologous intergenic regions
3. Definition, statistical analysis and plotting of pan- and core-genomes, LSEs ...
4. Comparative genome analyses, including analysis of the pan-genome structure (core, soft-core, shell and cloud) and identification of lineage-specific clusters



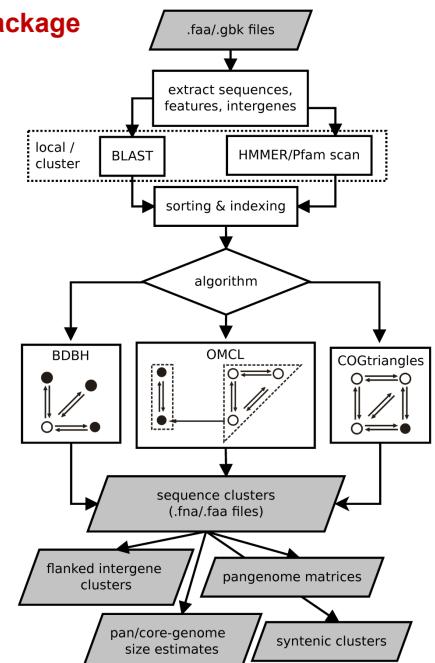
Overview of the GET_HOMOLOGUES package

Block 1: Extracting features from GBK files

Block 2: Blasting, sorting & indexing

Block 3: clustering

Block 4: parsing, statistical analysis and graphical display of pan-genomes and core-genomes with auxiliary scripts



get_homologues: a software package for microbial comparative genomics.

get_homologues.pl installation

- Written in [Perl](#), tested in 32 and 64-bit Linux, Mac OS X systems
- Installed along with dependencies using [install.pl](#)
- Bundled with precompiled binary files (COGtriangles, OMCL and BLAST+)
- Required dependencies (Perl modules)

Bio::Seq, Bio::SeqIO,

- Optional software and database dependencies:

hmmscan (from the HMMER3 package) for Pfam domain scanning

Pfam HMM library and hmmpress for db formatting

BerkleyDB (perl module)

R

get_homologues: a software package for microbial comparative genomics. [Running get_homologues.pl](#)

```
* usage: .../get_homologues.pl [options]
... cont.
* Options that control clustering:
-D require equal Pfam domain composition
when defining similarity-based orthology
-S min %sequence identity in BLAST query/subj pairs
-N min BLAST neighborhood correlation PubMed=18475320
-b compile core-genome with minimum BLAST searches

Options that control clustering:
-t report sequence clusters including at least t taxa
-a report clusters of sequence features in GenBank files
instead of default 'CDS' GenBank features

-g report clusters of intergenic sequences flanked by ORFs
in addition to default 'CDS' clusters
-f filter by %length difference within clusters

-r reference proteome .faa/.gbk file

-e exclude clusters with inparalogues
-x allow sequences in multiple COG clusters
-F orthoMCL inflation value
```

(best with -m cluster or -n threads)

(range [1-100], default=1 [BDBH|OMCL])
(range [0,1], default=0 [BDBH|OMCL])
ignores -c [BDBH])

(default t=numberOfTaxa,
t=0 reports all clusters [OMCL|COGS])
requires -d and .gbk files,
example -a 'tRNA,rRNA',
NOTE: uses blastn instead of blastp,
ignores -g,-D)
requires -d and .gbk files)

(range [1-100], by default sequence
length is not checked)
By default takes file with
least sequences; with BDBH sets
first taxa to start adding genes)
By default inparalogues are
included)
By default sequences are allocated
to single clusters [COGS])
(range [1-5], default=1.5 [OMCL])

get_homologues: a software package for microbial comparative genomics. [Running get_homologues.pl](#)

Input data: whole genome GenBank or FASTA files

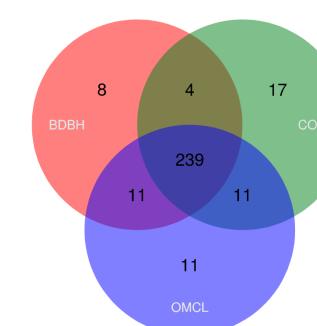
Typing \$./get_homologues.pl -h, on the terminal will show the basic options:

```
*usage: .../get_homologues.pl [options]
-h this message
-v print version, credits and checks installation
-d directory with input amino acid FASTA files (.faa) or (overrides -i)
GenBank files (.gbk), 1 per taxon; allows for new files to be added
there later, creates output folder named 'directory_homologues'
-i input amino acid FASTA file with [taxon names] in headers,
(required unless -d is set, creates output folder named 'file_homologues'
forces -m local)

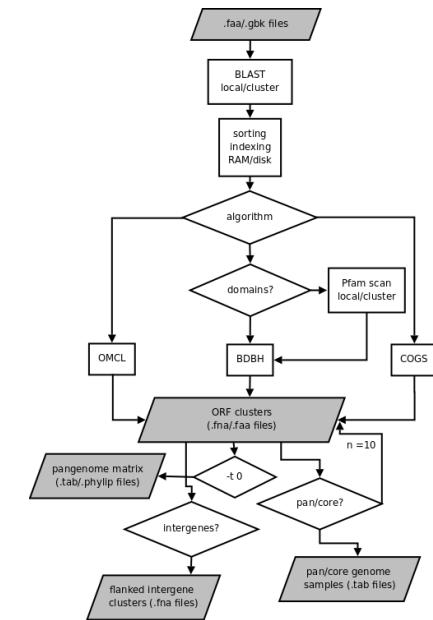
*Optional parameters:
-o only run BLAST/Pfam searches and exit (useful to pre-compute searches)
-c report genome composition analysis (follows order in -I file if enforced,
ignores -r,-t,-e)
-s save memory by using BerkeleyDB; default parsing stores sequence hits in RAM
-m runmode [local|cluster] (default local)
-I file with .faa/.gbk files in -d to be included (takes all by default,
requires -d)

*Algorithms instead of default bidirectional best-hits (BDBHs):
-G use COGtriangle algorithm (COGS, PubMed=20439257) (requires 3+ genomes|taxa)
-M use orthoMCL algorithm (OMCL, PubMed=12952885)
... output truncated
```

get_homologues: a software package for microbial comparative genomics. [Running get_homologues.pl – choosing a proper clustering algorithm](#)



Venn diagram showing the overlap between clusters of orthologous sequences produced by the three algorithms using the 119 *Stenotrophomonas* genomes analyzed by Vinuesa et al. 2018.



get_homologues: a software package for microbial comparative genomics.

Running get_homologues.pl – performing a core/pan-genome analysis

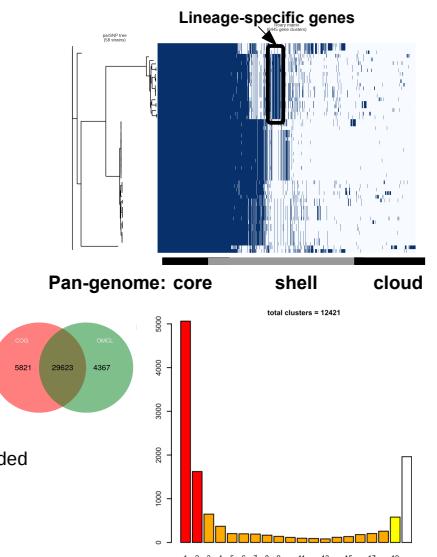
For a pangenome analysis use options:

-t 0 (get all clusters
of homologous sequences)

-M or -G (MCL or COGtriangles)

Note 1: the use of the BDBH clustering

Note 2: it is possible to compute a consensus pan-genomes using the intersection between the clusters detected by the COG and OMCL algorithms with the help of the [compare_clusters.pl](#) script, provided with the distribution.



learning about group-specific genes and lineage-specific gene expansions with parse_pangenome_matrix.pl (126 genes specifically found in EHEC and not in K2 - B group commensals)

genes pr

4631 predicted_hydrolase_faa
4638 predicted_reductase_faa
5036 predicted_antirepressor_protein_faa
5047_Sfp_A-_systemic_factor_protein_A-like_protein_faa
5081_T3SS_secured_effector_NleE_homolog_faa
5161_predicted_DNA_methylase_faa
5364_predicted_DNA-binding_protein_faa
5376_predicted_protein_faa
5419_predicted_terminase_large_subunit_faa
5672_predicted_receptor_protein_faa
5730_predicted_transmembrane_protein_faa
5783_predicted_terminase_small_subunit_faa
5785_predicted_head_protein-prohead_protease_faa
5788_predicted_DNA_packaging_protein_faa
5760_predicted_minor_tail_protein_faa
5765_predicted_tall_length_tape_measure_protein_faa
5775_T3SS_secured_effector_NleA-EspI_homolog_faa
6116_predicted_protein_faa
6148_predicted_phage_tail_tape_measure_protein_faa
6150_predicted_tall_protein_faa
6396_yeast_faa
7284_predicted_protein_faa
7284_predicted_4-hydroxybenzoate_decarboxylase_faa
7445_predicted_DNA_bind_protein_faa
7459_predicted_lipoprotein_faa
7496_T3SS_secured_effector_NleB_homolog_faa
7497_T3SS_secured_effector_NleE_homolog_faa
7498_EfaL_LtpA_protein_faa
7846_looser_peptidase_HopD_faa
7921_predicted_protein_faa

0053_lpFEE_faa
8187_escf_faa
8188_cedD2_faa
8189_escP_faa
8190_escP2_faa
8191_escP_faa
8192_escP_faa
8193_escD_faa
8194_escA_faa
8195_cecT_faa
8196_tir_faa
8197_mop_faa
8198_escf_faa

secreted_in_set_A_and_absent_in_B(126):	
dicted_hydrolases_faa	
dicted_reductase_faa	
dicted_repressor_protein_faa	
A_systemic_factor_protein_A-like_protein_faa	
S_secreted_effector_NieB_homolog_faa	
dicted_DnaJ_methylase_faa	
dicted_NADP-dependent_protein_faa	
dicted_protein_faa	
dicted_terminase_large_subunit_faa	
dicted_replcation_protein_faa	
dicted_protein_faa	
dicted_terminase_small_subunit_faa	
dicted_head_protein_protease_faa	
dicted_DNA_packaging_protein_faa	
dicted_minor_tail_protein_faa	
dicted_tall_length_tape_measure_protein_faa	
S_secreted_effector_NieA_EspI_homolog_faa	
dicted_protein_faa	
dicted_phage_tall_length_tape_measure_protein_faa	
dicted_tail_protein_faa	
dicted_protein_faa	
dicted_hexosemonosaccharidecarboxylase_faa	
dicted_DNA-bindin_protein_faa	
S_secreted_effector_NieB_homolog_faa	
S_secreted_effector_NieE_homolog_faa	
AA_protein_faa	
dicted_lipoprotein_faa	
dicted_protein_faa	
dicted_hemolysin_activating_lysine-acetyltransferase_HlyC_faa	
17707_hemolysin_A_faa	
18889_putative_permease_component_of_sugar_ABC_transport_system_faa	
18890_putative_permease_component_of_sugar_ABC_transport_system_faa	
19110_conserved_hypothetical_protein_faa	
19955_conserved_hypothetical_protein_faa	
29762_conserved_hypothetical_protein_faa	
32862_terW_faa	
32867_hypothetical_protein_faa	
32868_conserved_hypothetical_protein_faa	
32869_TerD_faa	
32870_terA_faa	
32871_terB_faa	
32872_terC_faa	
32873_terD_faa	
49048_Hypothetical_protein_faa	
60626_hypothetical_protein_faa	
64461_putative_helicase_faa	
67509_Mg_faa	
111250_T3SS_secreted_effector_EspI-X-h	
111252_T3SS_secreted_effector_EspN-h	
111697_horaine_attaching_effacing_associ	
112043_T3SS_secreted_effector_EspV-W_h	
112037_C4Axy_n-terminal_prolease_faa	
112035_putative_ArcA_family_regulatory_pro	
112382_hypothetical_protein_faa	
112791_putative_antirepressor_faa	
112820_hypothetical_protein_faa	
113505_T3SS_secreted_effector_NieG-hom	
115669_putative_tellurium_resistance_pro	
115671_hypothetical_protein_faa	
115690_putative_diacylglycerol_Kinase_faa	
115691_putative_membrane-anchored_pro	
115694_hypothetical_gamma-proteobact_UreF_faa	
115695_uridylate_alpha_subunit_UreF_faa	
115696_uridylate_alpha_subunit_UreF_faa	
115697_uridylate_accessory_protein_UreF_faa	
113698_uridylate_accessory_protein_UreF_faa	
113699_uridylate_accessory_protein_UreF_faa	
115021_hypothetical_protein_faa	
115088_putative_outer_membrane_precurs	
115141_putative_lambda_maystoyl_transf	
22062_head-tail_preambler_gib_P	
22144_putative_inner_membrane_protein_faa	
22155_protein_specific_repressor_faa	
224902_pofU_faa	
224911_putative_transcriptional_regulator_f	
224619_putative_inner_membrane_protein_faa	
222322_putative_ATP-binding_protein_faa	
22482_wab_faa	
223403_waiI_faa	

***E. coli* pangenome phylogeny estimated using the consensus panmatrix of presence/absence data for 46 fully sequenced genomes of *Escherichia* spp. / *Shigella* spp. and *Citrobacter* spp.**

Robust identification of orthologues and paralogues for microbial pan-genomics using GET_HOMOLOGUES: a case study of pIncA/C plasmids

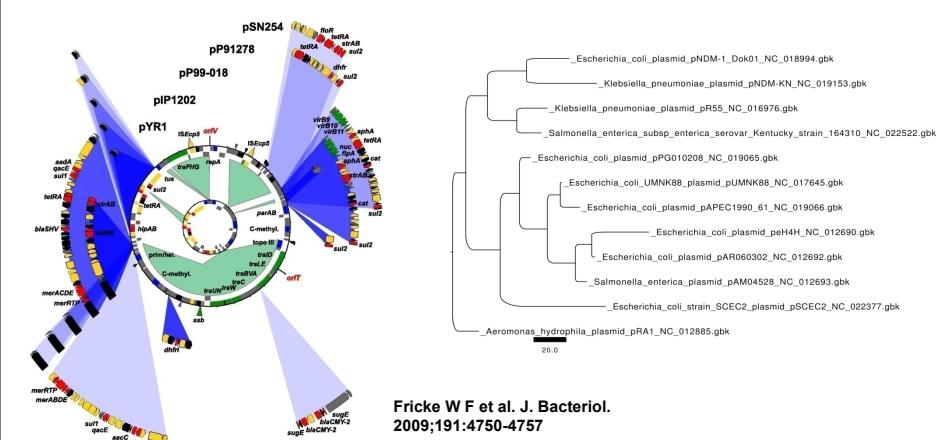
Pablo Vinyeta^{1*} and Bruno Contreras-Moreira^{1,2,3}

1 Centro de Ciencias Genómicas, Universidad Nacional Autónoma de México, Cuernavaca, Morelos, Mexico

2 Fundación ARAID, Zaragoza, Spain

3 Estación Experimental de Aula Dei, Consejo Superior de Investigaciones Científicas (EEAD-CSIC). Avda. Montaña, 1005, 50005 Zaragoza, Spain.

In: "Bacterial Pan-genomics". Methods in Molecular Biology. Marco Galardini, Alessio Mengoni and Marco Facci Ede, Humana Press, Springer, 2014. In press.



The GET_HOMOLOGUES tutorial:

1) cp the instructions script to your home

```
cp /home/vinuesa/genomica_UAEM/get_hom/code4_GET_HOMOLOGUES_OMICAS-UAEM-UNAM.txt .
```

Or download/visualize it from the GitHub repository:

https://github.com/vinuesa/OMICAS_UAEM/blob/master/docs/get_hom/code4_GET_HOMOLOGUES_OMICAS-UAEM-UNAM.txt

2) Open the script with nedit and keep a terminal open to issue the commands

```
gedit code4_GET_HOMOLOGUES_OMICAS-UAEM-UNAM.txt &
```

Analyses to be performed:

A pangenomic analysis of pIncA/C plasmids using GET_HOMOLOGUES

- Defining a robust core- and pan-genome of pIncA/C plasmids
- Exploring the gene space of pIncA/C plasmids: core, shell, cloud
- Pan-genome trees vs. core genome trees (supermatrices)
- Identifying lineage-specific genes in NDM-1 producing plasmids